

Mini-Project 6: Principal Component Analysis & Face Clustering

INSTRUCTOR: DANIEL L. PIMENTEL-ALARCÓN

DUE 04/23/2018

In this mini-project you will use principal component analysis (PCA) to cluster images of human faces under varying illuminations. To this end we will use the Yale dataset (`Yale_data.mat`), which includes the data array `Images` of size $48 \times 42 \times 64 \times 38$, containing 2432 photos of 38 subjects (64 photos per subject), each photo of size 48×42 . These images look like:



The main intuition is that the vectorized photos of the same subjects lie close to each other in the space of principal components, and so we will use a simple *nearest neighbor* approach on such space.

- Let $\mathbf{x}_{ij} \in \mathbb{R}^{2016}$ denote the j^{th} vectorized photo of the i^{th} subject, and let $\mathbf{X} \in \mathbb{R}^{2016 \times 2432}$ be the data matrix containing *all* the vectorized photos. Randomly split \mathbf{X} into *training* data $\mathbf{X}_a \in \mathbb{R}^{2016 \times 2204}$ and *testing* data $\mathbf{X}_b \in \mathbb{R}^{2016 \times 228}$, in such a way that \mathbf{X}_a contains 58 images of each subject, and \mathbf{X}_b contains 6 images of each subject.
- Let $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}$ denote the singular value decomposition of \mathbf{X}_a , such that $\mathbf{X}_a = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$.
- The j^{th} column of \mathbf{U} denotes the j^{th} principal vector. Display the 5 (unvectorized) leading principal vectors, often called *eigenfaces*.
- The diagonal entries in $\mathbf{\Sigma}$ denote the singular values. Plot their magnitudes. How many of them are significant?
- Let \mathbf{U}_r denote the matrix formed with the first r columns of \mathbf{U} , where r is your answer from (d). This way, \mathbf{U}_r spans the subspace containing *most* of the information of \mathbf{X} .
- Let $\mathbf{\Theta}_a \in \mathbb{R}^{r \times 2204}$ be the coefficient matrix of \mathbf{X}_a with respect to \mathbf{U}_r , such that $\mathbf{X}_a = \mathbf{U}_r \mathbf{\Theta}_a$, and similarly for $\mathbf{\Theta}_b \in \mathbb{R}^{r \times 228}$.

At this point we have transformed \mathbf{X}_a and \mathbf{X}_b into principal components space. More precisely, $\mathbf{\Theta}_a$ and $\mathbf{\Theta}_b$ are the representation of \mathbf{X}_a and \mathbf{X}_b with respect to the basis of principal vectors \mathbf{U}_r . Now we will classify

each column in \mathbf{X}_b (corresponding to a photo) by finding its nearest neighbor in \mathbf{X}_a in principal components space. That is, we will run nearest neighbor on Θ_b using Θ_a as reference.

- (h) Compute the *distance* matrix $\mathbf{D} \in \mathbb{R}^{2204 \times 228}$ between the columns of Θ_a and Θ_b . More precisely, the (i,j) th entry in \mathbf{D} contains the euclidian distance between the i th column of Θ_a and the j th column of Θ_b .
- (i) Display a few assignments (a photo in \mathbf{X}_b , and the photo in \mathbf{X}_a that was assigned to it). Does this seem to work?
- (j) Let $\hat{\mathbf{X}}$ be the projection of \mathbf{X} onto $\text{span}\{\mathbf{U}_r\}$. What is the normalized error $\|\mathbf{X} - \hat{\mathbf{X}}\|_F / \|\mathbf{X}\|_F$? What does this tell you?

I have created the following code to help you get started:

```

1   % © Daniel L. Pimentel-Alarcón, 2018, pimentel@gsu.edu
2   - close all; clear all; clc;
3
4   % ===== LOAD DATA =====
5   - load Yale_data.mat;
6   % FYI: Data array 'Images' contains 2432 photos of 38 subjects (64 photos
7   % per subject), each photo of size 48 x 42.
8
9   % ===== SPLIT BETWEEN TRAINING AND TESTING =====
10  - idx = randperm(64);
11  - Images_a = double(Images(:,:,idx(1:58),:));
12  - Images_b = double(Images(:,:,idx(59:64),:));
13
14  % ===== VECTORIZE DATA =====
15  % COMPLETE HERE: Vectorize data (both training and testing)
16
17  % ===== DO PCA OF TRAINING DATA =====
18  - [U,Sigma,V] = % COMPLETE HERE: Do SVD of training data
19
20  % ===== PLOT TOP EIGENFACES =====
21  - figure(1);
22  - for i=1:5,
23  -     subplot(1,5,i);
24  -     eigenface = % COMPLETE HERE: Unvectorize ith principal vector;
25  -     imagesc(eigenface);
26  -     colormap(gray);
27  -     axis image;
28  - end
29
30  % ===== PLOT SINGULAR VALUES =====
31  - figure(2);
32  % COMPLETE HERE: Plot singular values; consider using stem instead of plot
33
34  % ===== COMPUTE COEFFICIENTS IN PC SPACE =====
35  - r = % COMPLETE HERE: How many singular values are significant?
36  - U_r = U(:,1:r);
37  - Theta_a = % COMPLETE HERE: Compute coefficients of X_a
38  - Theta_b = % COMPLETE HERE: Compute coefficients of X_a
39

```

(continues below)

```

40 % ===== COMPUTE DISTANCE MATRIX OF COEFFICIENTS =====
41 - D = % COMPLETE HERE: Compute distance matrix
42 - figure(3);
43 - imagesc(D);
44 - colormap(gray);
45
46 % ===== NEAREST NEIGHBOR IN PC SPACE =====
47 - [~,closest] = % COMPLETE HERE: find closest training coefficient
48 % of each test coefficient
49
50 % ===== SEE A RANDOM RESULT =====
51 - choice = % COMPLETE HERE: Pick a test column randomly
52 - image_b = % COMPLETE HERE: Unvectorize test column
53 - image_a = % COMPLETE HERE: Unvectorize its corresponding training column
54
55 - figure(4);
56 - subplot(1,2,1);
57 - imagesc(image_b);
58 - colormap(gray);
59 - axis image;
60
61 - subplot(1,2,2);
62 - imagesc(image_a);
63 - colormap(gray);
64 - axis image;
65
66 % ===== PROJECTION ERROR =====
67 - X = [X_a X_b];
68 - Xhat = % COMPLETE HERE: Compute projection of X
69 - error = norm(X-Xhat,'fro')/norm(X,'fro')

```